

Regression Assignment

1. Problem statement:

- a. The problem statement is to compute the Insurance charges for a customer based on their **Age, BMI, No. of children, SEX, Smoker or Not a Smoker**

2. Understanding:

Independent variables: Age, BMI, No. of children, SEX, Smoker or Not a Smoker

Dependent: Charges

3. Total Number of rows: 1337

4. Total Number of columns: 6

5. Preprocessing:

- a. As SEX, SMOKER column has String (Nominal data), need to convert into 0/1, thus applied pandas library -> `get_dummies` method . Thus the dataset became:

	age	bmi	children	charges	sex_male	smoker_yes
0	19	27.900	0	16884.92400	False	True
1	18	33.770	1	1725.55230	True	False
2	28	33.000	3	4449.46200	True	False
3	33	22.705	0	21984.47061	True	False
4	32	28.880	0	3866.85520	True	False
...
1333	50	30.970	3	10600.54830	True	False
1334	18	31.920	0	2205.98080	False	False
1335	18	36.850	0	1629.83350	False	False
1336	21	25.800	0	2007.94500	False	False
1337	61	29.070	0	29141.36030	False	True

1338 rows × 6 columns

R2_Score:

As we have multiple inputs and output as numerical data, I have applied MLR, SVM, Decision Tree, Random Forest algorithms here:

Model Name	R ² _score
MLR	0.7894790349867009
SVM	0.6288792857320359
Decision Tree	0.8716934306026609
Random Forest	0.8736226875503907

Tuning the parameters:

SVM:

Kernel	C	Gamma	R ² _score
<i>Linear</i>	0.1	Scale	-0.080959968427891
<i>Poly</i>	0.1	Scale	-0.08830237655410711
<i>Rbf</i>	0.1	Scale	-0.08907451521042731
<i>Sigmoid</i>	0.1	scale	-0.08826991450485111
<i>Linear</i>	10	Scale	0.462468414233968
<i>Poly</i>	10	Scale	0.038716222760231345
<i>Rbf</i>	10	Scale	-0.03227329390671052
<i>Sigmoid</i>	10	scale	0.03930714378274347
<i>Linear</i>	100	Scale	0.6288792857320359
<i>Poly</i>	100	Scale	0.6179569624059795
<i>Rbf</i>	100	Scale	0.3200317832050831
<i>Sigmoid</i>	100	scale	0.5276103546510407
<i>Linear</i>	0.1	auto	-0.080959968427891
<i>Poly</i>	0.1	auto	-0.080959968427891
<i>Rbf</i>	0.1	auto	-0.08907451521042731
<i>Sigmoid</i>	0.1	auto	-0.08826991450485111
<i>Linear</i>	10	auto	0.462468414233968
<i>Poly</i>	10	auto	0.038716222760231345
<i>Rbf</i>	10	auto	-0.03227329390671052
<i>Sigmoid</i>	10	auto	0.03930714378274347
<i>Linear</i>	100	auto	0.6288792857320359
<i>Poly</i>	100	auto	0.6179569624059795
<i>Rbf</i>	100	auto	0.3200317832050831
<i>Sigmoid</i>	100	auto	0.5276103546510407

DecisionTree:

Criterion	Max_depth	Splitter	R ² _score
<i>squared_error</i>	5	<i>best</i>	0.8420667783377725
<i>friedman_mse</i>	5	<i>best</i>	0.8247592981719993
<i>absolute_error</i>	5	<i>best</i>	0.8549334830978529
<i>poisson</i>	5	<i>best</i>	0.8370067462772068
<i>squared_error</i>	50	<i>best</i>	0.6923136340310512
<i>friedman_mse</i>	50	<i>best</i>	0.6930369835887366
<i>absolute_error</i>	50	<i>best</i>	0.6862874546706454
<i>Poisson</i>	50	<i>best</i>	0.719693906627027
<i>squared_error</i>	5	random	0.8491041436762746
<i>friedman_mse</i>	5	random	0.8006347068373014
<i>absolute_error</i>	5	random	0.867493203036949
<i>poisson</i>	5	random	0.8716934306026609
<i>squared_error</i>	50	random	0.7367163929864847
<i>friedman_mse</i>	50	random	0.6971743652804676
<i>absolute_error</i>	50	random	0.7005039657893241
<i>Poisson</i>	50	random	0.6910698158676895

RandomForest:

Criterion	n_estimators	Max_features	R ² _score
<i>squared_error</i>	50	sqrt	0.870029408972246
<i>absolute_error</i>	50	sqrt	0.8716413352284358
<i>friedman_mse</i>	50	sqrt	0.868283465836858
<i>poisson</i>	50	sqrt	0.8733655509389101
<i>squared_error</i>	50	Log2	0.8621383503359108
<i>absolute_error</i>	50	Log2	0.8736226875503907
<i>friedman_mse</i>	50	Log2	0.8629236318262976
<i>poisson</i>	50	Log2	0.8715022263954666
<i>squared_error</i>	100	sqrt	0.8682671249870114
<i>absolute_error</i>	100	sqrt	0.873183586394964
<i>friedman_mse</i>	100	sqrt	0.8678844813076738
<i>poisson</i>	100	sqrt	0.8694685934948518
<i>squared_error</i>	100	Log2	0.8694726783046982
<i>absolute_error</i>	100	Log2	0.8708500790901526
<i>friedman_mse</i>	100	Log2	0.868583614709796

<i>poisson</i>	100	Log2	0.870665321846255
----------------	-----	------	-------------------

Algorithm chosen:

The model built using **Decision Tree** algorithm has the best R^2 _score when compared to other algorithms.